

Reservoir computing for spatiotemporal signal classification without trained output weights *

Ashley Prater

Air Force Research Laboratory, Information Directorate, Rome NY USA †

July 20, 2016

Abstract

Reservoir computing is a recently introduced machine learning paradigm that has been shown to be well-suited for the processing of spatiotemporal data. Rather than training the network node connections and weights via backpropagation in traditional recurrent neural networks, reservoirs instead have fixed connections and weights among the ‘hidden layer’ nodes, and traditionally only the weights to the output layer of neurons are trained using linear regression. We claim that for signal classification tasks one may forgo the weight training step entirely and instead use a simple supervised clustering method based upon principal components of norms of reservoir states. The proposed method is mathematically analyzed and explored through numerical experiments on real-world data. The examples demonstrate that the proposed may outperform the traditional trained output weight approach in terms of classification accuracy and sensitivity to reservoir parameters.

1 Introduction

Reservoir computing is a recently developed bio-inspired machine learning paradigm for the processing of spatiotemporal data [1, 2]. In the language of neural networks, a reservoir is collection of hidden layer nodes with nonlinear recurrent dynamics, where the nodes are sparsely connected with fixed weights that are not trained to fit specific data. Because the weights are fixed, using a reservoir requires only a simple initialization step, as opposed to more traditional recurrent neural networks whose weights and connections must be learned in a tedious backpropagation training step [3]. The property of fixing the reservoir connections and weights has many benefits, including ease of initialization, along with having the ability to quickly adapt to new data and applications.

Reservoirs, like all recurrent neural networks, are based on the premise that the state of the reservoir at a particular time should depend on the current value of the input signal, along with recent inputs and reservoir states. To be an effective method for computation, a reservoir should map input data into a sufficiently high-dimensional space. It is desirable for a reservoir to operate ‘at the edge of chaos’ [4], so dissimilar inputs are sufficiently separated in the reservoir node states, yet inputs with only small perturbation-like differences do not stray too far apart. Reservoir dynamics demonstrate long short-term memory [5], so any individual point-wise errors in a signal will not corrupt the entire reservoir response.

Two types of reservoirs that have emerged in literature include echo state networks (ESNs) and time-delay reservoirs (TDRs). An ESN uses randomly, yet sparsely, connected nodes with randomly assigned fixed weights [1, 2, 6]. A TDR uses a cyclic topology, where each node provides data to exactly one other node, and has fixed, non-random weights [7, 8, 9]. The output layer of both ESN and TDR-type reservoirs traditionally use linear output weights, trained on a labeled dataset using least squares or ridge regression [1, 10, 11]. This method has an easy training phase, and is computationally cheap to use in the testing phase. However, it can be sensitive to reservoir parameters and dataset characteristics and prone to overfitting. If the training

*This work was cleared for public release by Wright Patterson Air Force Base Public Affairs on 11 Apr 2016. Case Number: 88ABW-2016-1812.

†ashley.prater.3@us.af.mil

dataset has large intra-class variation, or if the classes are not well-separated, then it may be difficult to find a collection of weights to discriminate the classes well.

In this research, a simple supervised clustering method based on principal components is proposed for use in classification tasks using ESNs and TDRs. The method used is based upon comparing the norm of a reservoir response of a test signal against the principal components of the norms of reservoir states for classes of labeled training data. The clustering method has slightly higher computational complexity than using trained output weights to classify new input signals, however it may achieve higher classification accuracy while being less sensitive to reservoir type, size, and feedback strength. We present a rigorous analysis of the clustering method, including two theorems characterizing the upper bound of the difference in reservoir responses for two input signals, with the upper bound in terms of the input signals, the reservoir type, and the user-generated parameters. Moreover, we explore the difference in performance of the two methods through numerical simulations performed using a real-world dataset for both ESNs and TDRs for various reservoir parameters. In every simulation, the clustering approach outperforms the trained output weights in terms of both accuracy and CPU time required to classify test signals.

The following notation is used in this work. For a collection of signals $\{u\}$, the j^{th} element in the collection is denoted by $u^{(j)}$. Training sets are partitioned into K classes. Let \mathcal{C}_k be the collection of indices of signals in the k^{th} class. That is, $u^{(j)}$ is in the k^{th} class iff $j \in \mathcal{C}_k$. For a vector v , the ℓ_2 norm is given by $\|v\|_2 = (\sum_j v_j^2)^{1/2}$. For a matrix A , $\rho(A)$ is the spectral radius, i.e. the largest absolute value of an eigenvalue of A . We use $\mathcal{O}(\cdot)$ with the standard ‘big O’ meaning, that $f(x) = \mathcal{O}(g(x))$ if there exists $M > 0$ and $x_0 \in \mathbb{R}$ such that $|f(x)| \leq M|g(x)|$ for all $x \geq x_0$.

2 Reservoir Computing Models for Classification

Suppose $u \in \mathbb{R}^T$ is an input signal of length T , possibly after the application of a multiplexing mask, and say $u(t)$ is the value of u at time t . The values of the N reservoir nodes at time t are called the *reservoir states* and are denoted by the vectors $X(t) \in \mathbb{R}^N$, one vector for each t . The n^{th} entry of these vectors, $X_n(t)$ denote the state of the n^{th} reservoir node at time t . The dynamics of the ESN and TDR architectures are described by the following models:

$$\text{ESN: } X(t) = f(W_{\text{in}}u(t) + W_{\text{res}}X(t-1)) \quad (1)$$

$$\text{TDR: } X_n(t) = \begin{cases} f(\alpha u(t) + \beta X_{N-1}(t-1)), & \text{if } n = 0 \\ X_{n-1}(t-1), & \text{if } n \in \{1, 2, \dots, N-1\} \end{cases} \quad (2)$$

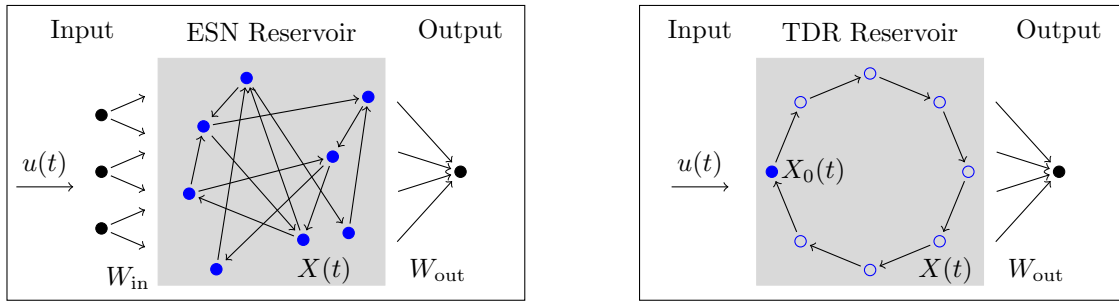


Figure 1: Representations of two architectural variants of reservoirs with output weights, the echo state network (left) and time delay reservoir (right).

For ease of notation, suppose each reservoir type has N nodes. In the ESN topology, the vector $W_{\text{in}} \in \mathbb{R}^N$ weights the input signal feeding into the nodes, while the matrix $W_{\text{res}} \in \mathbb{R}^{N \times N}$ determines the fixed connections and weights among the nodes. That is, node m feeds into node n weighted by the $(n, m)^{\text{th}}$ entry of W_{res} in the ESN model. The TDR has $N-1$ virtual nodes, corresponding to $n = 1, 2, \dots, N-1$, and one physical node for $n = 0$. The parameter α is the input gain, and β is the attenuation value. Notice in the TDR the node values are simply passed along the reservoir unchanged except at the physical node. Models of the ESN and TDR reservoirs are shown in Figure 1.

The function f in Equations (1) and (2) is a nonlinear activation function. Typical choices for f include sinusoidal, logistic, sigmoidal, and piecewise linear functions.

Time-multiplexing the inputs is a common preprocessing step in time-delay reservoir systems [7, 8, 9, 12, 13]. The multiplexing mask is applied as follows. Suppose the raw inputs are z_1, z_2, z_3, \dots , and consider a mask m of length L . Then the multiplexed input u is defined via $u((k-1)L+1 : kL) = z_k m$, that is each raw input z_k is multiplied by the vector m and concatenated to form the multiplexed input. The purpose of the multiplexing mask in TDRs is several-fold. A non-constant mask helps to increase the dimensionality of the reservoir, yielding richer dynamics [9, 13]. Furthermore, since the inputs are all passed to the reservoir only via the head node, the mask allows several virtual nodes to process values from a single raw input vector at once, as random ESNs do by design [12, 13]. An additional benefit is that it helps to ‘slow down’ TDRs, many of which are physically implemented as optical devices that would process the raw data much faster than one can sample the outputs [8, 9, 14].

Two approaches for interpreting the reservoir outputs for supervised classification tasks are described in the subsections below. The first describes the traditional approach using trained output weights, and the second describes a method of clustering the reservoir node states.

2.1 Trained linear output weights

A classical method to interpret the results of a reservoir is to train a collection of output weight matrices $W_{\text{out}}(t) \in \mathbb{R}^{K \times N}$ at each time t of interest [1, 10, 11, 15] that map reservoir states close to an appropriate ‘indicator’ vector. That is, choose a collection of times of interest $\Omega \subseteq \{1, 2, \dots, T\}$, and let $X^{(j)}(t) \in \mathbb{R}^M$ denote the reservoir nodes at time t driven by the j^{th} element in the training dataset using either Equation (1) or (2). Then the collection of output weights at time t is

$$W_{\text{out}}(t) = \underset{W \in \mathbb{R}^{K \times N}}{\text{argmin}} \left\{ \sum_{j \in \text{Tr}} \|d_j - W X^{(j)}(t)\|_2^2 + \lambda \|W\|_2^2 \right\}, \quad (3)$$

where each $d_j(t) \in \mathbb{R}^K$ is an indicator vector, having all zeros entries except for a 1 in the k^{th} position if $j \in \mathcal{C}_k$.

Equipped with the collection of output weights, a test pattern with reservoir node states $X(t) \in \mathbb{R}^N$ is determined to belong to the k^{th} class if the K -vector

$$D = \sum_{t \in \Omega} \omega_t W_{\text{out}}(t) X(t) \quad (4)$$

has maximal element in the k^{th} row. Typically the classification weights ω_t are chosen to be 1.

Algorithm 1 (To classify a signal using trained linear output weights)

Initialization: Input the fixed parameters λ, Ω .

Training: Generate the vectors $X^{(j)}(t) \in \mathbb{R}^N$ for each j and $t \in \Omega$, using Equation (1) or (2), then find the collections $\{W_{\text{out}}(t) : t \in \Omega\}$ as in Equation (3).

Testing: Let $u \in \mathbb{R}^T$ be a new test pattern.

1. Compute the corresponding reservoir nodes $\{X(t) \in \mathbb{R}^N : t \in \Omega\}$ using Equation (1) or (2).
 2. Compute the vector D as in Equation (4).
 3. Say u is in the k^{th} class if $D(k) \geq D(\ell)$ for all indices $\ell \in \{1 : K\}$.
-

The computational cost of determining the class of a new pattern using the ‘Testing’ phase of Algorithm 1 is determined as follows. Assume that the matrices $W_{\text{out}}(t)$ are given, and do not include its derivation in the cost evaluation. To drive the reservoir and find the nodes $X(t)$ of a new test pattern requires $\mathcal{O}(N^2 T)$ multiplications using the ESN dynamics in Equation (1), or $\mathcal{O}(T)$ multiplications using the TDR dynamics in Equation (2). Although only the reservoir node values at times $t \in \Omega$ are of interest, one must drive the reservoir using the full set of times. To find the vector D requires $\mathcal{O}(KN|\Omega|)$ multiplications, and

finally $\mathcal{O}(K)$ comparisons are needed to determine the maximal element. Overall, this leads to a complexity of $\mathcal{O}(N^2T + KN|\Omega|)$ when using ESN-type reservoirs and a complexity of $\mathcal{O}(T + KN|\Omega|)$ when using TDR-type reservoirs.

2.2 Classification via Clustering with Principal Components

The underlying idea for the training method (3) is that similar inputs to the reservoir produce similar outputs, even after the non-linear high-dimensional processing is applied. Under this assumption, it is feasible that one could classify data using a clustering method without the use of the trained output weights. Therefore, we propose the following method using the principal components of norms of reservoir responses to perform classification. Let \mathcal{C}_k be the collection of indices of training patterns that belong to the k^{th} class. Find the reservoir states $X^{(j)}(t) \in \mathbb{R}^N$ for all $j \in \mathcal{C}_k$, $t \in \Omega$, and k , and compute the vectors $b_j \in \mathbb{R}^{|\Omega|}$, where

$$b_j(i) = \left\| X^{(j)}(t_i) \right\|_2^2, \quad t_i \in \Omega. \quad (5)$$

For each class, concatenate the vectors b_j to form matrices $B_k \in \mathbb{R}^{|\Omega| \times |\mathcal{C}_k|}$. Since the input training patterns belong to the same class, the columns of each B_k should exhibit similar characteristics. Suppose $U_k \in \mathbb{R}^{|\Omega| \times R}$ is the matrix of the first R principal components of B_k . For any new test patterns with corresponding vector b , one can say that the pattern belongs to the k^{th} class if U_k describes b well, i.e. if

$$\|(I - U_k U_k^*) b\| \leq \|(I - U_\ell U_\ell^*) b\|, \quad \forall \ell.$$

Algorithm 2 (To classify a signal using clustering via principal components)

Initialization: Input the collection of times Ω and the number of principal components to consider R .
Training:

1. Generate the vectors $X^{(j)}(t) \in \mathbb{R}^N$ for all $j \in \mathcal{C}_k$, $t \in \Omega$, and each k using either Equation (1) or (2).
2. Compute the matrices B_k with columns as in Equation (5).
3. Compute U_k , the first R principal components of B_k .

Testing: Let $u \in \mathbb{R}^T$ be a new test pattern.

1. Generate the reservoir states $X(t) \in \mathbb{R}^N$ for each t according to Equation (1) or (2).
 2. Compute the vector $b \in \mathbb{R}^{|\Omega|}$ according to Equation (5).
 3. For each k , compute $d_k = \|(I_{|\Omega|} - U_k U_k^*) b\|_2^2$.
 4. Say u belongs to the k^{th} class if $d_k \leq k_\ell$ for all ℓ .
-

The computational cost of determining the class of a new pattern using the ‘Testing’ phase of Algorithm 2 is determined as follows. Assume that the matrices $I - U_k U_k^* \in \mathbb{R}^{|\Omega| \times |\Omega|}$ are precomputed during the training phase. As in Algorithm 1, the cost to drive the reservoir and find the nodes $X(t)$ requires $\mathcal{O}(N^2T)$ multiplications using the ESN, or $\mathcal{O}(T)$ multiplications using the TDR. To compute the vector b in Step 2 requires $\mathcal{O}(N|\Omega|)$ multiplications, and to compute the values $\{d_k\}$ in Step 3 requires $\mathcal{O}(K|\Omega|^2)$ multiplications. Finally, to determine the class of u in Step 4 requires K comparisons. Overall, this leads to a complexity of $\mathcal{O}(N^2T + N|\Omega| + K|\Omega|^2)$ for ESN-type reservoirs, and a complexity of $\mathcal{O}(T + N|\Omega| + K|\Omega|^2)$ for TDR-type reservoirs. Since the parameters can vary in magnitude, the dominant term in the complexity depends on the particular set-up used.

3 Analysis of Reservoir Behavior

The clustering method proposed in Algorithm 2 will be more accurate if small variations in the input signals lead to bounded differences in reservoir states, while large discrepancies in inputs are mapped farther apart. To confidently use this approach, we must characterize reservoir responses for similar inputs.

Several studies of reservoir performance based on the type of reservoir architecture, chosen parameters, as well as the characteristics of the input data have been performed, with evidence that some combinations of the aforementioned factors can seriously degrade performance [2, 7, 10, 16]. However, the metrics used in the reservoir computing literature tend to be only experimentally investigated. To explore how well the reservoir response separates classes, the separation ratio [6, 17], point-wise separation [2, 18], and class separation [19] have been used. These all measure how well a reservoir can separate inputs from distinct classes, by having distances between disparate classes large while keeping similar inputs close. Similarly, to measure how effectively a reservoir can process a particular dataset, researchers use the universal approximation property [2] kernel quality [19, 20, 21], reservoir capacity [22], and the Echo State Property [1]. These measures and properties concern the representation of inputs within the reservoir response and the reconstructability of an input signal from reservoir states. For robustness to noise, generalization rank [21] or the Lyapunov coefficient [17, 19, 20, 23, 24, 25] are considered.

Although the reservoir dynamics (1) and (2) have simple descriptions, rigorous treatment of their behavior have proven difficult, with few results so far. In Proposition 3 of [1], the distance between two reservoir states at a given time is bounded in terms of the reservoir states at the previous timestep and the spectral radius of the reservoir weights. Although mathematically proven, this Proposition covers only randomly connected ESNs incrementing one timestep with activation functions of the form $f(x) = \tanh(x)$. Theorem 3.5 of [26] bounds the distance between two output vectors of a TDR, determined using linear read-out weights, in terms of the reservoir parameters and the behavior of the input signals. In the Theorems below, we prove upper bounds for distances between two reservoir responses to different inputs in terms of reservoir parameters and the behavior of the inputs for both ESNs and TDRs, and in more generality than the results given in [1] and [26].

For readability, let us first introduce some notation. Let $u^{(j)}(t)$ denote the j^{th} input at time t , with corresponding reservoir states $X^{(j)}(t)$. Let $\delta_{i,j,t} = |u^{(i)}(t) - u^{(j)}(t)|$ be the difference between two input signals at time t , and let $\varepsilon_{i,j,t} = \|x^{(i)}(t) - x^{(j)}(t)\|$ be the distance between the corresponding node states at time t . Suppose $\bar{\delta}_{i,j} = \sup\{\delta_{i,j,t} : t \in \mathbb{R}\}$ is bounded for each pair (i, j) , and that the nonlinear activation function f is Lipschitz continuous with optimal Lipschitz constant L . Finally, let $[\cdot]_n$ denote a vector whose entries run over the range of the variable n .

Theorem 1. *Suppose the reservoir node states are determined using the ESN dynamics from Equation (1). If $\rho(W_{\text{res}})$ is the spectral radius of W_{res} , then the distance between the reservoir nodes at time t corresponding to two input signals $u^{(i)}$ and $u^{(j)}$ satisfies*

$$\varepsilon_{i,j,t} \leq L\bar{\delta}_{i,j}\|W_{\text{in}}\| \frac{1 - (L\rho(W_{\text{res}}))^{t+1}}{1 - L\rho(W_{\text{res}})}.$$

Proof. By Equation (1) and the Lipschitz continuity of f ,

$$\begin{aligned} \varepsilon_{i,j,t} &= \|X^{(i)}(t) - X^{(j)}(t)\| \\ &= \left\| f\left(W_{\text{in}}u^{(i)}(t) + W_{\text{res}}X^{(i)}(t-1)\right) - f\left(W_{\text{in}}u^{(j)}(t) + W_{\text{res}}X^{(j)}(t-1)\right) \right\| \\ &\leq L \left\| W_{\text{in}}[u^{(i)}(t) - u^{(j)}(t)] + W_{\text{res}}[X^{(i)}(t-1) - X^{(j)}(t-1)] \right\| \\ &\leq L\|W_{\text{in}}\|\delta_{i,j,t} + L\rho(W_{\text{res}})\varepsilon_{i,j,t-1}. \end{aligned}$$

Since $\varepsilon_{i,j,-1} = 0$, it follows by induction that

$$\varepsilon_{i,j,t} \leq L\|W_{\text{in}}\| \sum_{r=0}^t (L\rho(W_{\text{res}}))^r \delta_{i,j,t-r} \leq L\bar{\delta}_{i,j}\|W_{\text{in}}\| \frac{1 - (L\rho(W_{\text{res}}))^{t+1}}{1 - L\rho(W_{\text{res}})}. \quad \square$$

Theorem 2. *Suppose the reservoir node states are determined using the TDR dynamics from Equation (2). Then the distance between the reservoir nodes corresponding to two input signals $u^{(i)}$ and $u^{(j)}$ satisfies*

$$\varepsilon_{i,j,t} \leq \alpha\bar{\delta}_{i,j}L\sqrt{N} \frac{1 - (\beta L)^{\lfloor t/N \rfloor + 1}}{1 - \beta L}.$$

Proof. By Equation (2) and the Lipschitz continuity of f ,

$$\begin{aligned}
\varepsilon_{i,j,t} &= \left\| [X_n^{(i)}(t)]_n - [X_n^{(j)}(t)]_n \right\| \\
&= \left\| [X_0^{(i)}(t-n)]_n - [X_0^{(j)}(t-n)]_n \right\| \\
&\leq \alpha L \left\| [u^{(i)}(t-n) - u^{(j)}(t-n)]_n \right\| + \beta L \left\| [X_{N-1}^{(i)}(t-n-1) - X_{N-1}^{(j)}(t-n-1)]_n \right\| \\
&\leq \alpha L \bar{\delta}_{i,j} \sqrt{N} + \beta L \left\| [X_n(t-N)]_n \right\| \\
&= \alpha L \bar{\delta}_{i,j} + \beta L \varepsilon_{i,j,t-N}.
\end{aligned}$$

Let $r \in \{0, 1, \dots, N-1\}$ be the remainder when t is divided by N . By induction on the inequality above,

$$\varepsilon_{i,j,t} \leq (\beta L)^{\lfloor t/N \rfloor} \varepsilon_{i,j,r} + \alpha \bar{\delta}_{i,j} L \sum_{k=0}^{\lfloor t/N \rfloor - 1} (\beta L)^k.$$

Since $r < N$, the n^{th} reservoir node at time r can be characterized by

$$X_n(r) = \begin{cases} f(\alpha u(r-n)), & \text{if } n \leq r \\ 0, & \text{if } n > r \end{cases},$$

yielding $\varepsilon_{i,j,r} \leq \alpha \bar{\delta}_{i,j} L \sqrt{r+1} \leq \alpha \bar{\delta}_{i,j} L \sqrt{N}$. Therefore

$$\varepsilon_{i,j,t} \leq \alpha \bar{\delta}_{i,j} L \sqrt{N} \sum_{k=0}^{\lfloor t/N \rfloor} (\beta L)^k = \alpha \bar{\delta}_{i,j} L \sqrt{N} \frac{1 - (\beta L)^{\lfloor t/N \rfloor + 1}}{1 - \beta L}. \quad \square$$

Theorems 1 and 2 show that for input signals with small pointwise discrepancies and well-chosen reservoir parameters, their associated reservoir state norms cluster well. However, the Theorems do not guarantee that very distinct inputs are mapped to dissimilar reservoir node state norms. For this, we turn to the separation ratio, introduced in [6] and further explored in [17]. For completeness, we include it here, modified for both Algorithm 1 and Algorithm 2. For Algorithm 1, operator on the reservoir responses themselves, and for Algorithm 2 consider the norms of the reservoir responses.

Define the center of mass of the reservoir states of the K^{th} class in the training set at time t as $M_k(t)$,

$$M_k(t) = \begin{cases} \frac{1}{|\mathcal{C}_k|} \sum_{j \in \mathcal{C}_k} X^{(j)}(t), & \text{for Algorithm 1,} \\ \frac{1}{|\mathcal{C}_k|} \sum_{j \in \mathcal{C}_k} \|X^{(j)}(t)\|, & \text{for Algorithm 2.} \end{cases}$$

The inter-class distance is computed the same for both algorithms. It is defined as the average distance between pairs of class means at each time step:

$$d(t) = \frac{1}{K^2} \sum_{k=1}^K \sum_{\ell=1}^K \|M_k(t) - M_\ell(t)\|,$$

The intra-class variance is the average variance within each class at each time step:

$$v(t) = \begin{cases} \frac{1}{K} \sum_{k=1}^K \frac{1}{|\mathcal{C}_k|} \sum_{j \in \mathcal{C}_k} \|M_k(t) - X^{(j)}(t)\|, & \text{for Algorithm 1,} \\ \frac{1}{K} \sum_{k=1}^K \frac{1}{|\mathcal{C}_k|} \sum_{j \in \mathcal{C}_k} |M_k(t) - \|X^{(j)}(t)\||, & \text{for Algorithm 2.} \end{cases}$$

Then the separation ratio at time t is defined as

$$\text{Sep}(t) = \frac{d(t)}{1 + w(t)}. \quad (6)$$

The larger $\text{Sep}(t)$ is, the better the separation among the classes at time t .



Figure 2: A subset of images from the USPS handwritten digit dataset.

4 Example

Handwritten digits are classified using the trained linear output weights in Algorithm 1 and the using the principal components method in Algorithm 2. The data used are from the United States Postal Service (USPS) database, obtained from [27]. A sample of these images is shown in Figure 2. Each image in the dataset is a 16×16 pixel 8-bit grayscale image, reshaped as a 256 length column vector taking values in $[0, 1]$. The data are split in ten classes of 1100 images each, representing the digits 0 through 9. For each simulation presented below, 400 images in each class are randomly selected to form the training set, while the remaining 7000 images (700 from each class) are used as the test set. Although the nearby pixel behavior is not preserved in the horizontal direction by transforming each image into a column vector, the correlations are still present in the reservoir response due to the long short term memory property.

All experiments are implemented in MATLAB R2013a on a node with 2 Intel Xeon 5650 CPUs with 8 cores at 2.67 GHz with 8GB RAM.

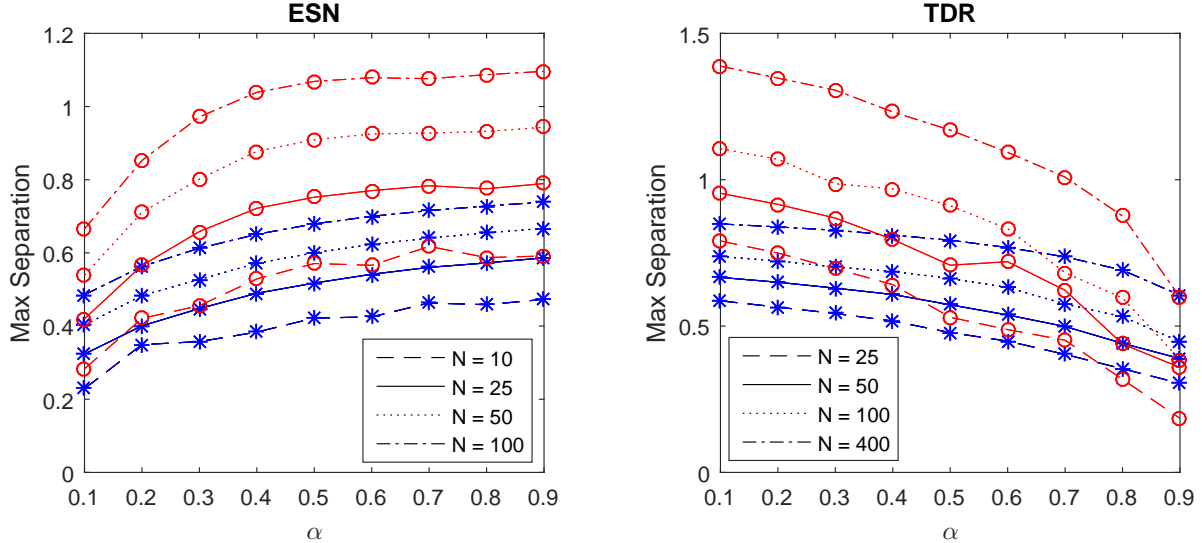


Figure 3: The maximum separation ratio from Equation (6) attained by the reservoirs on sample training sets, adapted for Algorithm 1 in blue ‘*’ and for Algorithm 2 in red ‘o’, for several values of α and N . Results for the ESN-style reservoirs are in the left plot, and results for TDR-style reservoirs are in the right plot.

4.1 Experiment Setup

The ESN reservoirs are set up with N nodes, for $N \in \{10, 25, 50, 100\}$. The input weights are $W_{\text{in}} = [\alpha \ \alpha \ \dots \ \alpha]^\top \in \mathbb{R}^N$, where α ranges over the set $\{0.1, 0.2, \dots, 0.9\}$. The reservoir weights $W_{\text{res}} \in \mathbb{R}^{N \times N}$ are randomly chosen with 20% density, and scaled so that the largest eigenvalue is $0.9999(1 - \alpha)$. No mask is used with ESN reservoirs, so $T = 256$, $\Omega = \{1, 2, \dots, 256\}$, and $K = 10$.

The TDR-type reservoirs are set up with N nodes, for $N \in \{25, 50, 100, 400\}$. Again, the parameter α appearing in Equation (2) ranges over the set $\{0.1, 0.2, \dots, 0.9\}$. The inputs are multiplexed with a mask of length $N - 1$ randomly taking values from $\{\pm 1\}$, but the reservoir is sampled only every $N - 1$ time-steps. Therefore $T = 256(N - 1)$ and $\Omega = \{r(N - 1) + 1 : r = 0, 1, \dots, 255\}$ with $|\Omega| = 256$.

For each simulation, 400 images from each class are randomly selected to form the training dataset, however, the same training dataset selection is used for each pair (N, α) . The nonlinear activation function is chosen to be $f(x) = \sin(x)$ throughout. For the trained linear output weights, the regularization parameters $\lambda = 10^{-4}$ and $\lambda = 10^{-10}$ are used.

4.2 Results

The results of these simulations are presented in Figures 3-5.

Figure 3 plots the maximum separation ratio from Equation (6) attained by ESN and TDR reservoirs for both Algorithms for the selections of parameters α and N . Notice that the reservoirs do not separate this data particularly well, but the norms of reservoir responses that are used in Algorithm 2 tend to be slightly better separated than the vector responses used in Algorithm 1.

The top two plots of Figure 4 show the average classification accuracy, and the bottom two plots give the time required to classify all 7000 images in the test set using both reservoir types with both Algorithm 1 and Algorithm 2 for the parameters α, N and λ . The results for the clustering approach presented in Algorithm 2 are denoted by red ‘ \circ ’. The results for the trained output weights using Algorithm 1 use blue ‘ $*$ ’ (for $\lambda = 10^{-4}$) or cyan ‘ \square ’ (for $\lambda = 10^{-10}$).

The clustering approach always achieves a higher classification accuracy than the trained linear output weights, but takes only about 35-40 seconds longer to classify all 7000 images. Notice the clustering approach is fairly robust to the choice of reservoir and parameters N and α . The trained linear output weights are more sensitive to N and α , and are inversely related to the separation ratio given in Figure 3.

The computational complexity of the two Algorithms can be seen in the ‘Time’ plots of Figure 4. For ESNs, both Algorithms have a quadratic dependency on N , but Algorithm 2 takes a bit longer also having a quadratic dependency on $|\Omega| = T$. For TDRs, the linear dependence on N for both Algorithms is evident in the plot. The TDR has a longer runtime than the ESN since a mask is used with the TDR, increasing T by a factor of $N - 1$.

The clustering approach in Algorithm 2 was also applied to the raw input dataset without using a reservoir. Over 100 trials, the average accuracy of the clustering method applied to the raw input data is 95.27%, which is smaller than the average accuracy attained by Algorithm 2 using an ESN or TDR. This suggests that the clustering method is well-suited to this problem, but processing the data in a reservoir improves accuracy for most parameter choices since the reservoir preserves the spatial correlations well.

Figure 5 displays the inequalities presented in Theorems 1 and 2, measuring the discrepancy of reservoir activations at time t for similar inputs. The two input signals were randomly selected from the class of ‘3s’. The values shown in the figure are the found by dividing out the right hand side of the inequality, giving

$$\varepsilon_{i,j,t} / \left(L \bar{\delta}_{i,j} \|W_{\text{in}}\| \frac{1 - (L\rho(W_{\text{res}}))^{t+1}}{1 - L\rho(W_{\text{res}})} \right)$$

in the left image, and

$$\varepsilon_{i,j,t} / \left(\alpha \bar{\delta}_{i,j} L \sqrt{N} \frac{1 - (\beta L)^{\lfloor t/N \rfloor + 1}}{1 - \beta L} \right)$$

in the right image, both plotted against t . The inequalities in the theorems are clearly satisfied since they are well below 1, however the upper limits could be further refined in future research.

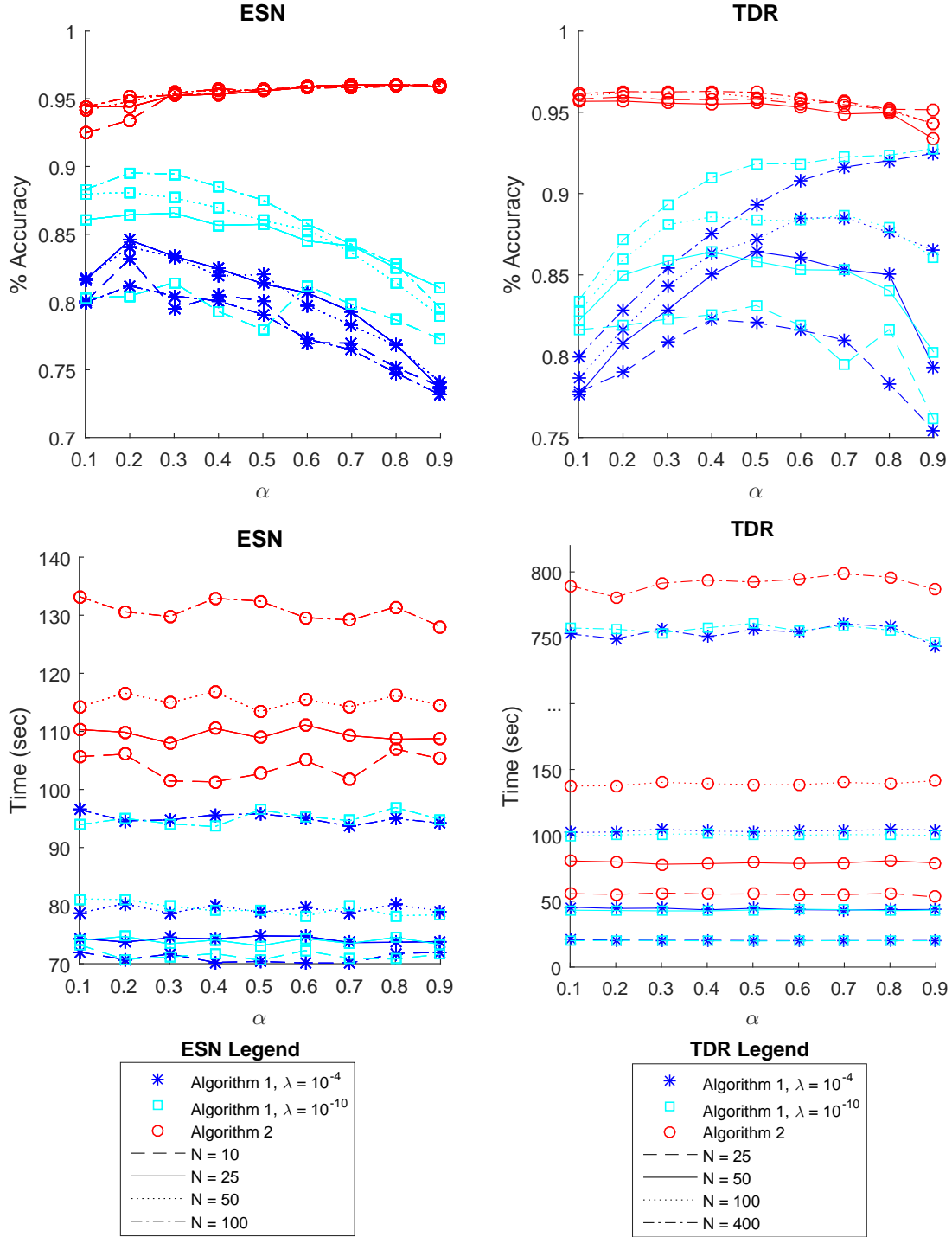


Figure 4: A comparison of the performance of the proposed method using Algorithm 2 and the method of trained linear output weights using Algorithm 1 on ESN and TDR reservoirs for various parameters α and reservoir size N . The top two figures present the classification accuracy on the test set, and the bottom two figures present the total time required in seconds to classify the entire test set of 7000 images. The plots with red '○' denote results from Algorithm 2, and the plots with blue '*' or cyan '□' denote results from Algorithm 1 for different regularization parameters.

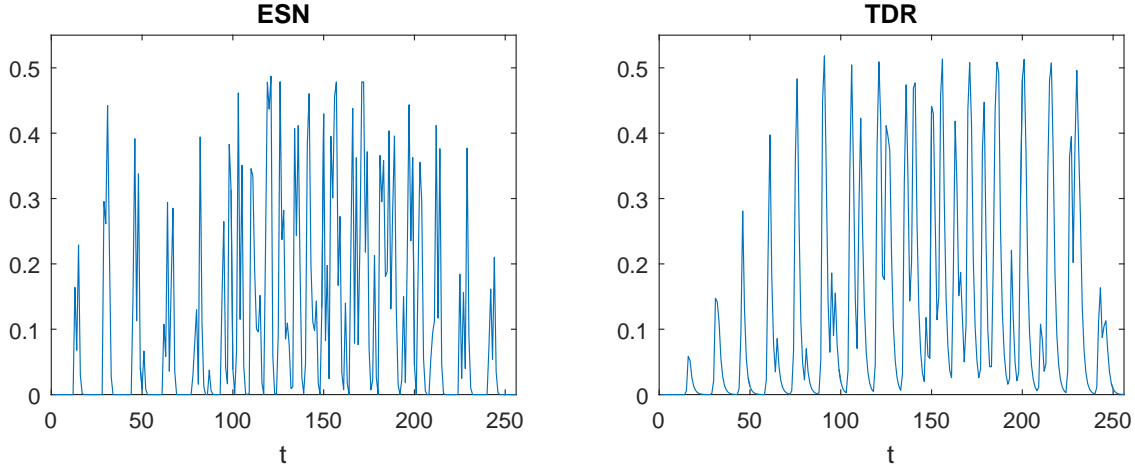


Figure 5: Ratios of the inequalities from Theorems 1 and 2 for randomly selected i and j , plotted against values of t , with $N = 100$ and $\alpha = 0.5$.

5 Conclusion

This work theoretically and experimentally explored a method to classify spatiotemporal patterns using the principal components of norms of reservoir states on a training set. The proposed method was compared to the traditional method using trained linear output weights for two types of reservoir topologies using several parameter selections. In the numerical experiments, the proposed method achieved better classification accuracy on the test set, but took a bit longer to complete computations. The proposed method loses some information since it considers norms of reservoir state vectors, but this leads to more robustness with respect to reservoir type and size, as well as parameter choice.

A basic implementation of both methods was used so the fundamental principles could be compared. More sophisticated implementations could be used in future work, and may improve speed and accuracy for both methods. These adaptations could include selecting better training sets, introducing subclasses to reduce intra-class variation and improve class separation, using optimally designed masks for TDRs [13], refining the reservoir connections and weights [28, 29], improving selection of parameters (spectral radius, reservoir size, feedback strength, regularization parameter) and subsequent solving of trained output weights.

Acknowledgements

This research was supported by Air Force Office of Scientific Research [LRIR:15RICOR122].

Any opinions, findings and conclusions or recommendations expressed in this material are those of the author and do not necessarily reflect the view of the United States Air Force.

References

- [1] H. Jaeger, The ‘echo state’ approach to analysing and training recurrent neural networks - with an erratum note, Tech. Rep. GMD Report Number 148, Fraunhofer Institute for Autonomous Intelligent Systems (2011).
- [2] W. Maass, H. Natschlager, H. Markram, Real-time computing without stable states: a new framework for neural computation based on perturbations, *Neural Comput.* 14 (2002) 2531–2560. doi:10.1162/089976602760407955.
- [3] P. Werbos, Backpropagation through time: What it does and how to do it, *Proc. IEEE* 78 (1990) 1550–1560. doi:10.1109/5.58337.

- [4] N. Bertschinger, H. Natschlager, Real-time computation at the edge of chaos in recurrent neural networks, *Neural Comput.* 6 (2004) 1413–1436. doi:10.1162/089976604323057443.
- [5] H. Jaeger, Long short-term memory in echo state networks: Details of a simulation study, Tech. Rep. 27, Jacobs University Bremen (2012).
- [6] E. Goodman, D. Ventura, Spatiotemporal pattern recognition via liquid state machines, in: 2006 International Joint Conference on Neural Networks (IJCNN), 2006, pp. 3848–3853. doi:10.1109/ijcnn.2006.246880.
- [7] L. Grigoryeva, J. Henriques, L. Larger, J. Ortega, Optimal nonlinear information processing capacity in delay-based reservoir computers, *Sci. Rep.* 5 (12858). doi:10.1038/srep12858.
- [8] S. Ortín, M. C. Soriano, L. Pesquera, D. Brunner, D. San-Martin, I. Fischer, C. R. Mirasso, J. M. Gutierrez, A unified framework for reservoir computing and extreme learning machines based on a single time-delayed neuron, *Sci. Rep.* 5. doi:10.1038/srep14945.
- [9] Y. Paquot, F. Duport, A. Smerieri, J. Dambre, B. Schrauwen, M. Haelterman, S. Massar, Optoelectronic reservoir computing, *Sci. Rep.* 2. doi:10.1038/srep00287.
- [10] A. Goudarzi, P. Banda, M. Lakin, C. Teuscher, D. Stefanovic, A comparative study of reservoir computing for temporal signal processing, Tech. rep., University of New Mexico (2014). URL [arXiv:1401.2224](https://arxiv.org/abs/1401.2224)
- [11] M. Lukoševičius, H. Jaeger, Reservoir computing approaches to recurrent neural network training, *Comp. Sci. Rev.* 3 (2009) 127–149.
- [12] L. Appletant, Reservoir computing based on delay-dynamical systems, Ph.D. thesis, Vrije Universiteit Brussel, Universitat de les Illes Balears (May 2012).
- [13] L. Appletant, G. Van der Sande, J. Danckaert, I. Fischer, Constructing optimized binary masks for reservoir computing with delay systems, *Sci. Rep.* 4 (3629). doi:10.1038/srep03629.
- [14] F. Duport, A. Smerieri, A. Akrouf, M. Haelterman, S. Massar, Virtual optical reservoir computing, *Advanced Photonics, OSA Technical Digest (JM5A.40)*. doi:10.1364/BGPP.2014.JM5A.40.
- [15] M. Lukoševičius, A practical guide to applying echo state networks, in: G. Montavon, et al. (Eds.), *NN: Tricks of the Trade*, 2nd Edition, Springer-Verlag Berlin Heidelberg, 2012, pp. 650–686.
- [16] Y. Paquot, J. Dambre, B. Schrauwen, M. Haelterman, S. Massar, Reservoir computing: A photonic neural network for information processing, in: *Proc. SPIE Nonlinear Optics and Applications IV*, 2010. doi:10.1117/12.854050.
- [17] T. Gibbons, Unifying quality metrics for reservoir networks, in: 2010 International Joint Conference on Neural Networks (IJCNN), 2010, pp. 1–7. doi:10.1109/IJCNN.2010.5596307.
- [18] B. Schrauwen, D. Verstraeten, J. Van Campenhout, An overview of reservoir computing: Theory, applications and implementations, in: *ESANN 2007 proceedings - European Symposium on Artificial Neural Networks*, Bruges Belgium, 2007, pp. 25–27.
- [19] J. Chrol-Cannon, Y. Jin, On the correlation between reservoir metrics and performance for time series classification under the influence of synaptic plasticity, *PLOS ONE* 9. doi:10.1371/journal.pone.0101792.
- [20] R. Legenstein, W. Maass, Edge of chaos and prediction of computational performance for neural circuit models, *Neural Networks* 20 (2007) 323–334. doi:10.1016/j.neunet.2007.04.017.
- [21] M. Soriano, D. Brunner, M. Escalona-Morán, C. Mirasso, I. Fischer, Minimal approach to neuro-inspired information processing, *Front. Comput. Neurosci.* 9. doi:10.3389/fncom.2015.00068.

- [22] J. Dambre, D. Verstraeten, B. Schrauwen, S. Massar, Information processing capacity of dynamical systems, *Sci. Rep.* 2. doi:10.1038/srep00514.
- [23] D. Verstraeten, B. Schrauwen, D. M., D. Stroobandt, An experimental unification of reservoir computing methods, *Neural Networks* 20 (2007) 391–403. doi:10.1016/j.neunet.2007.04.003.
- [24] B. Schrauwen, L. Buesing, R. Legenstein, On computational power and the order-chaos phase transition in reservoir computing, in: *Proc. of NIPS 2008, Advances in Neural Information Processing Systems*, 2009.
- [25] B. Nils, H. Natschlager, Real-time computation at the edge of chaos in recurrent neural networks, *Neural Comput.* 16 (2004) 1413–1436. doi:10.1162/089976604323057443.
- [26] C. DiMarco, Reservoir computing dynamics for single nonlinear node with delay line structure, *Tech. rep.* (2015).
URL [arXiv:1510.03800](https://arxiv.org/abs/1510.03800)
- [27] S. Roweis, Data for MATLAB hackers., www.cs.nyu.edu/roweis/data.html, accessed: 16 May 2014.
- [28] H. Jaeger, M. Lukosevicius, D. Popovici, U. Siewert, Optimization and applications of echo state networks with leaky-integrator neurons, *Neural Networks* 20 (3) (2007) 335–352.
- [29] D. Norton, D. Ventura, Improving liquid state machines through iterative refinement of the reservoir, *Neurocomputing* 73 (2010) 2893–2904.